

Compressive Recovery of Sparse Precision Matrices

Etienne Lasalle

with



Titouan Vayer



Rémi Gribonval



Paulo Gonçalves

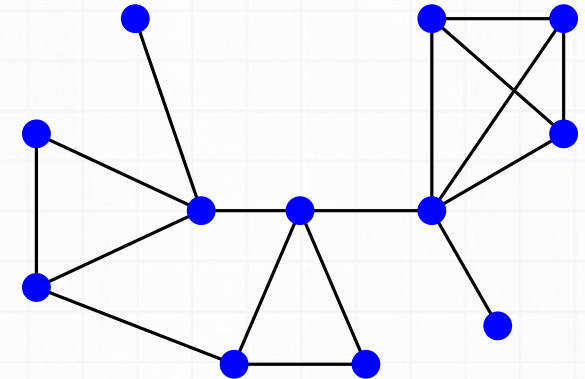
Journées de Statistiques, Bordeaux.
May 2024.

The problem



Data

Find the interactions
between the d variables



Θ : precision matrix

(Conditional dependencies)

$$\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma = \Theta^{-1})$$

Gaussian model

Framework:

- high dimension : d is large
- Sparsity in Θ

GLASSO approach [FHT08], [BGA08]

Naive estimator of Θ (MLE): $\hat{\Sigma}^{-1}$ (Not sparse)

Graphical Lasso

i.e., ℓ_1 -penalized maximum likelihood estimator

$$\hat{\Theta}_{\text{GL}} \triangleq \arg \min_{\Theta \succ 0} \left\{ -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle + \lambda \|\Theta\|_{1,\text{off}} \right\},$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$$

$$\|\Theta\|_{1,\text{off}} = \sum_{i < j} |\Theta_{ij}|$$

λ : regularization parameter

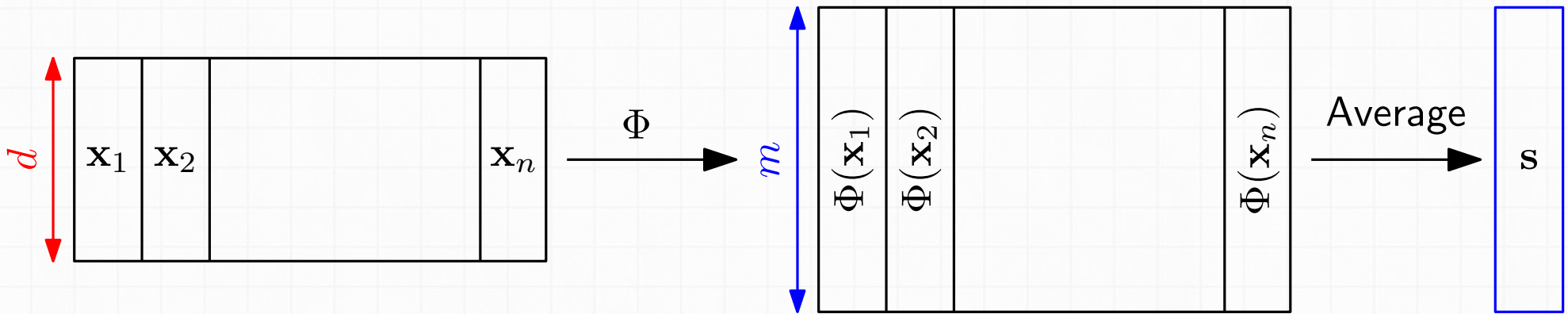
Memory cost : $\mathcal{O}(d^2)$ (storage of $\hat{\Sigma}$).

Computational cost : $\mathcal{O}(d^3)$.

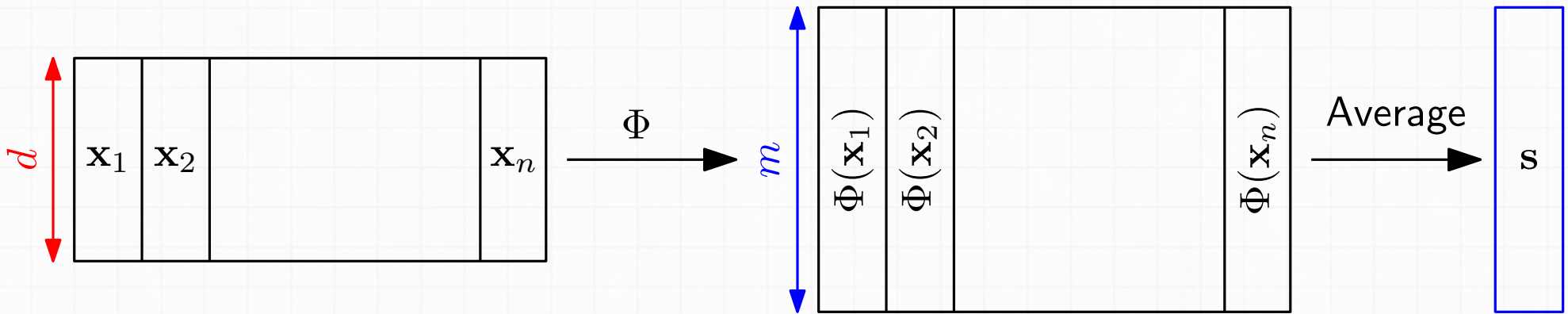
[FHT08] Friedman et al., *Sparse inverse covariance estimation with the graphical lasso*, (2008)

[BGA08] Banerjee et al., *Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data*, (2008)

Sketching approach (compression) [GCK+21]

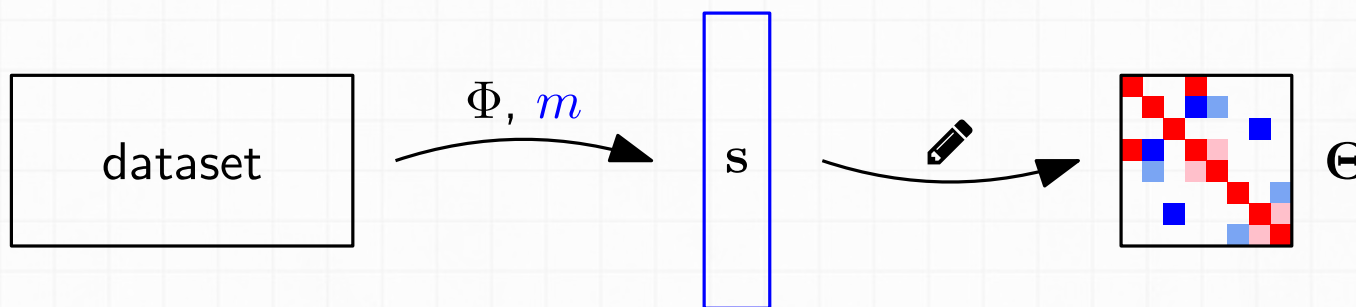


Sketching approach (compression) [GCK+21]

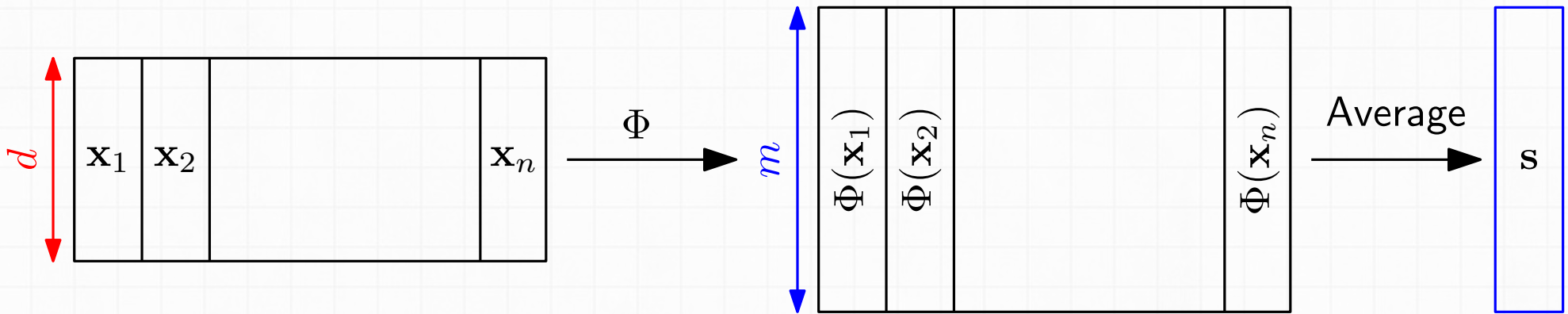


Questions :

- How to define s ? (Ideally : $\|\Theta\|_0 \lesssim m \ll d^2$)
- Sketch : is it efficient (computation and memory)?
- Practical decoder?
- Is it efficient?

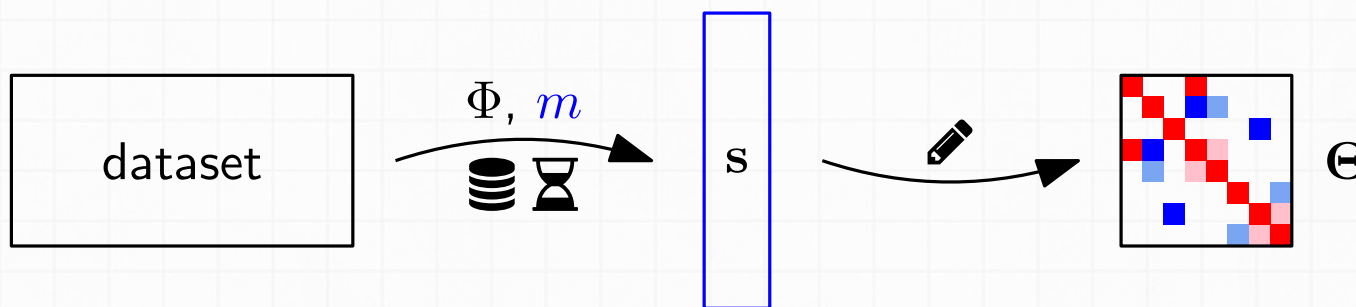


Sketching approach (compression) [GCK+21]

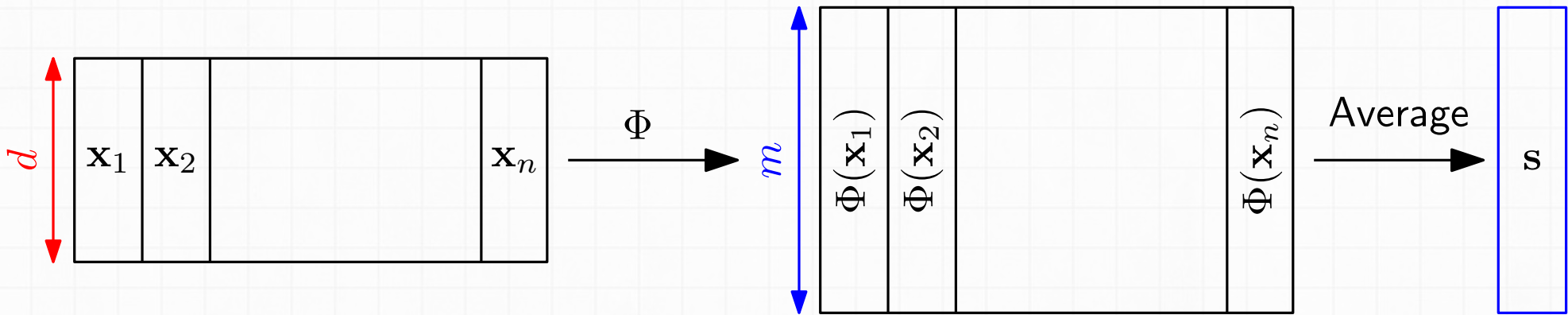


Questions :

- How to define s ? (Ideally : $\|\Theta\|_0 \lesssim m \ll d^2$)
- Sketch : is it efficient (computation and memory)?
- Practical decoder?
- Is it efficient?

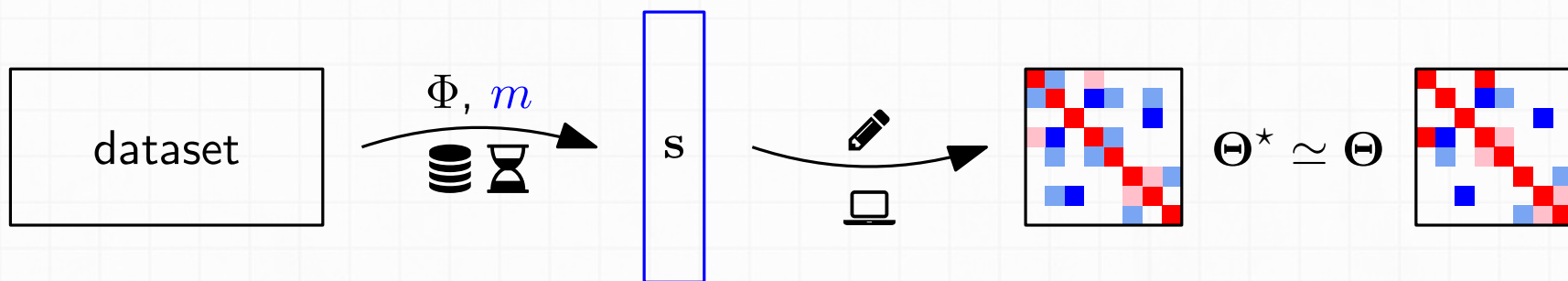


Sketching approach (compression) [GCK+21]

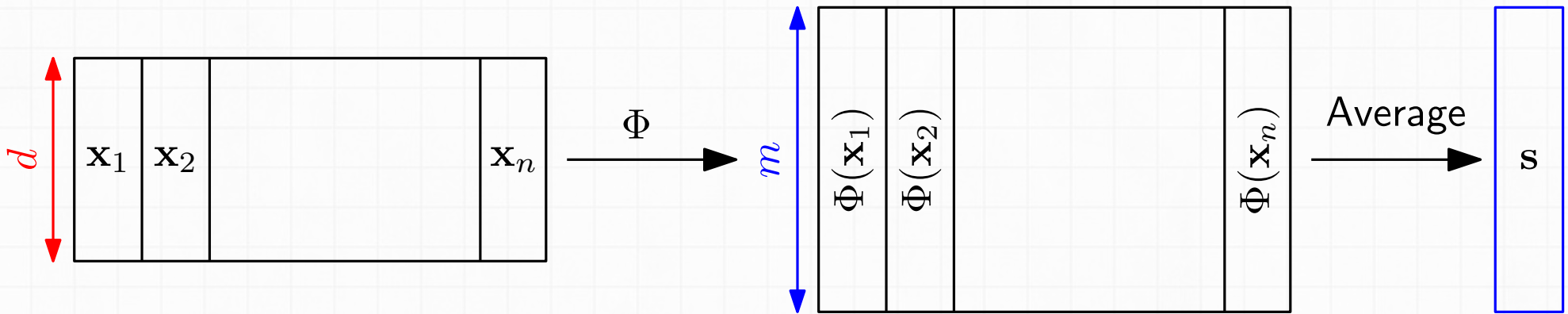


Questions :

- How to define s ? (Ideally : $\|\Theta\|_0 \lesssim m \ll d^2$)
- Sketch : is it efficient (computation and memory)?
- Practical decoder?
- Is it efficient?

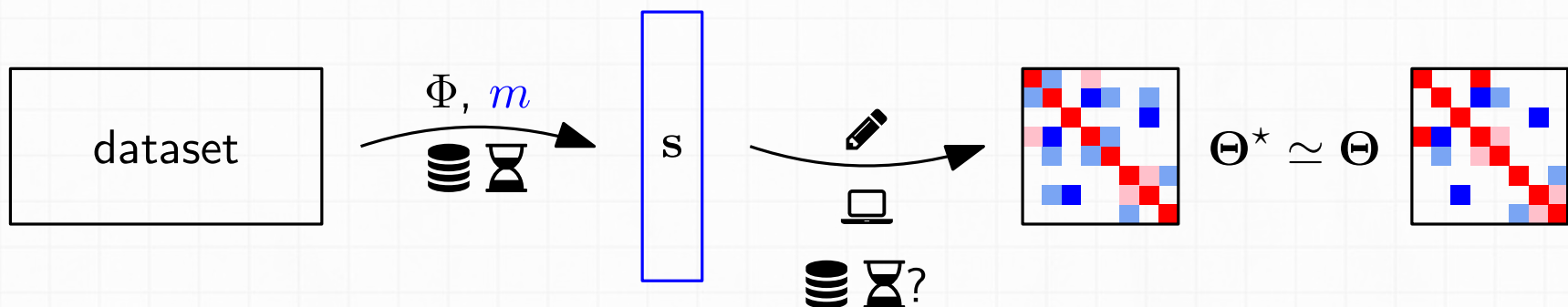


Sketching approach (compression) [GCK+21]



Questions :

- How to define s ? (Ideally : $\|\Theta\|_0 \lesssim m \ll d^2$)
- Sketch : is it efficient (computation and memory)?
- Practical decoder?
- Is it efficient?



Random Rank-One Projections (ROP)

$$\Phi(\mathbf{x}) \triangleq m^{-1} (\langle \mathbf{A}_1, \mathbf{x}\mathbf{x}^\top \rangle, \dots, \langle \mathbf{A}_m, \mathbf{x}\mathbf{x}^\top \rangle)^\top$$

with $\mathbf{A}_1, \dots, \mathbf{A}_m$ random $d \times d$ matrices.

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i) = m^{-1} (\langle \mathbf{A}_j, \hat{\Sigma} \rangle)_j = \mathcal{A}(\hat{\Sigma})$$

Noisy linear measurement of Σ .
(Compressed Sensing, Inv. Pb)

Random Rank-One Projections (ROP)

$$\Phi(\mathbf{x}) \triangleq m^{-1} (\langle \mathbf{A}_1, \mathbf{x}\mathbf{x}^\top \rangle, \dots, \langle \mathbf{A}_m, \mathbf{x}\mathbf{x}^\top \rangle)^\top$$

with $\mathbf{A}_1, \dots, \mathbf{A}_m$ random $d \times d$ matrices.

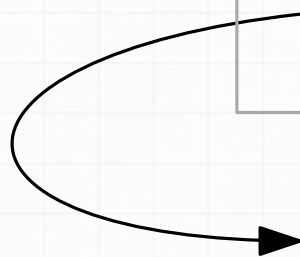
$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i) = m^{-1} (\langle \mathbf{A}_j, \hat{\Sigma} \rangle)_j = \mathcal{A}(\hat{\Sigma})$$

Noisy linear measurement of Σ .
(Compressed Sensing, Inv. Pb)

Computational and memory cost of $\Phi(x)$:

	Comp. ⌚	Memory 🗄️
\mathbf{A}_j full rank	$\mathcal{O}(md^2)$	$\mathcal{O}(md^2)$
\mathbf{A}_j rank one $\mathbf{A}_j = \mathbf{a}_j \mathbf{a}_j^\top$	$\mathcal{O}(md)$	$\mathcal{O}(md)$

[CZ15]
[CCG15]



$$\Phi(\mathbf{x}) = m^{-1} (|\mathbf{a}_1^\top \mathbf{x}|^2, \dots, |\mathbf{a}_m^\top \mathbf{x}|^2)^\top$$

[CZ15] Cai and Zhang, *Rop: Matrix recovery via rank-one projections*, (2015)

[CCG15] Chen et al., *Exact and stable covariance estimation from quadratic sampling via convex programming*, (2015)

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i) = m^{-1} (\langle \mathbf{A}_j, \hat{\Sigma} \rangle)_j = \mathcal{A}(\hat{\Sigma})$$

Definition: Restricted Isometry Property (RIP).

The *sketching operator* \mathcal{A} satisfies $\text{RIP}(\delta, \mathfrak{S})$, if for every $\underline{\Sigma}_1, \underline{\Sigma}_2 \in \mathfrak{S} \subseteq S_d$,

$$(1 - \delta) \|\underline{\Sigma}_1 - \underline{\Sigma}_2\| \leq \|\mathcal{A}(\underline{\Sigma}_1) - \mathcal{A}(\underline{\Sigma}_2)\| \leq (1 + \delta) \|\underline{\Sigma}_1 - \underline{\Sigma}_2\|.$$

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i) = m^{-1} (\langle \mathbf{A}_j, \hat{\Sigma} \rangle)_j = \mathcal{A}(\hat{\Sigma})$$

Definition: Restricted Isometry Property (RIP).

The *sketching operator* \mathcal{A} satisfies $\text{RIP}(\delta, \mathfrak{S})$, if for every $\Sigma_1, \Sigma_2 \in \mathfrak{S} \subseteq S_d$,

$$(1 - \delta) \|\Sigma_1 - \Sigma_2\| \leq \|\mathcal{A}(\Sigma_1) - \mathcal{A}(\Sigma_2)\| \leq (1 + \delta) \|\Sigma_1 - \Sigma_2\|.$$

Consequence :

“Optimal Decoder” : $\Sigma^* \triangleq \arg \min_{\tilde{\Sigma} \in \mathfrak{S}} \|\mathcal{A}(\tilde{\Sigma}) - \mathbf{s}\|$

$$\mathfrak{S}_{k,a,b} \triangleq \{ \Sigma \succ 0, \Theta = \Sigma^{-1} \text{ is } (d + 2k)\text{-sparse, } \text{spec}(\Theta) \subseteq [a, b] \}.$$

$$\|\mathcal{A}(\mathbf{M})\| = \|\mathcal{A}(\mathbf{M})\|_1 \text{ and } \|\mathbf{M}\| \text{ defined from } \mathcal{L}(\mathbf{a}) \text{ s.t. } \mathbb{E} [\|\mathcal{A}(\mathbf{M})\|_1] = \|\mathbf{M}\|.$$

$$\|(\Sigma^*)^{-1} - \Theta\|_{\text{Fro}} \leq C b^2 d \|\Sigma^* - \Sigma\| \leq \frac{2Cb^2 d}{(1-\delta)} \|\mathcal{A}(\hat{\Sigma}) - \mathcal{A}(\Sigma)\|_1.$$

Theoretical recovery guarantees

$\mathfrak{S}_{k,a,b} \triangleq \{ \Sigma \succ 0, \Theta = \Sigma^{-1} \text{ is } (d + 2k)\text{-sparse, } \text{spec}(\Theta) \subseteq [a, b] \}$.

$\|\mathcal{A}(\mathbf{M})\| = \|\mathcal{A}(\mathbf{M})\|_1$ and $\|\mathbf{M}\|$ defined from $\mathcal{L}(\mathbf{a})$ s.t. $\mathbb{E} [\|\mathcal{A}(\mathbf{M})\|_1] = \|\mathbf{M}\|$.

Theorem : [VLGG23]

Let \mathcal{A} be our sketching operator and $\mathbf{a}_1, \dots, \mathbf{a}_m \stackrel{i.i.d.}{\sim}$ Gaussian or Unif. on S_d .
 $\forall \delta \in]0, 1[, \exists C = C(\delta, b/a)$ s.t. whenever

$$m \geq C(d + 2k) \log d,$$

\mathcal{A} satisfies $\text{RIP}(\delta, \mathfrak{S}_{k,a,b})$ with high probability.

Theoretical recovery guarantees

$\mathfrak{S}_{k,a,b} \triangleq \{ \Sigma \succ 0, \Theta = \Sigma^{-1} \text{ is } (d + 2k)\text{-sparse, } \text{spec}(\Theta) \subseteq [a, b] \} .$

$\|\mathcal{A}(\mathbf{M})\| = \|\mathcal{A}(\mathbf{M})\|_1$ and $\|\mathbf{M}\|$ defined from $\mathcal{L}(\mathbf{a})$ s.t. $\mathbb{E} [\|\mathcal{A}(\mathbf{M})\|_1] = \|\mathbf{M}\|$.

Theorem : [VLGG23]

Let \mathcal{A} be our sketching operator and $\mathbf{a}_1, \dots, \mathbf{a}_m \stackrel{i.i.d.}{\sim}$ Gaussian or Unif. on S_d .
 $\forall \delta \in]0, 1[, \exists C = C(\delta, b/a)$ s.t. whenever

$$m \geq C(d + 2k) \log d,$$

\mathcal{A} satisfies $\text{RIP}(\delta, \mathfrak{S}_{k,a,b})$ with high probability.

$$\mathfrak{S}_{k,a,b} \longleftrightarrow \mathfrak{S}_{k,\kappa_0} \triangleq \{ \Sigma \succ 0, \Theta = \Sigma^{-1} \text{ is } (d + 2k)\text{-sparse, } \kappa(\Theta) \leq \kappa_0 \},$$
$$C(\delta, b/a) \longleftrightarrow C(\delta, \kappa_0)$$

Practical limitations

Encoding

$$m \geq C(d + 2k) \log d$$

	Comp. ⌚	Memory 🗄️
\mathbf{A}_j full rank	$\mathcal{O}(md^2)$	$\mathcal{O}(md^2)$
\mathbf{A}_j rank one	$\mathcal{O}(md)$	$\mathcal{O}(md)$
$\mathbf{A}_j = \mathbf{a}_j \mathbf{a}_j^\top$	$\mathcal{O}(d^2 \log d)$	$\mathcal{O}(d^2 \log d)$

Decoding

“Optimal Decoder”:

$$\Sigma^* \triangleq \arg \min_{\tilde{\Sigma} \in \mathcal{S}} \|\mathcal{A}(\tilde{\Sigma}) - \mathbf{s}\|$$

(Unsolvable)

Practical limitations

Encoding

$$m \geq C(d + 2k) \log d$$

	Comp. ⌚	Memory 🗄️
\mathbf{A}_j full rank	$\mathcal{O}(md^2)$	$\mathcal{O}(md^2)$
\mathbf{A}_j rank one $\mathbf{A}_j = \mathbf{a}_j \mathbf{a}_j^\top$	$\mathcal{O}(md)$ = $\mathcal{O}(d^2 \log d)$	$\mathcal{O}(md)$ = $\mathcal{O}(d^2 \log d)$
Structured \mathbf{A}_j	$\mathcal{O}(m \log d)$	$\mathcal{O}(m)$

Decoding

“Optimal Decoder”:

$$\Sigma^* \triangleq \arg \min_{\tilde{\Sigma} \in \mathcal{S}} \|\mathcal{A}(\tilde{\Sigma}) - \mathbf{s}\|$$

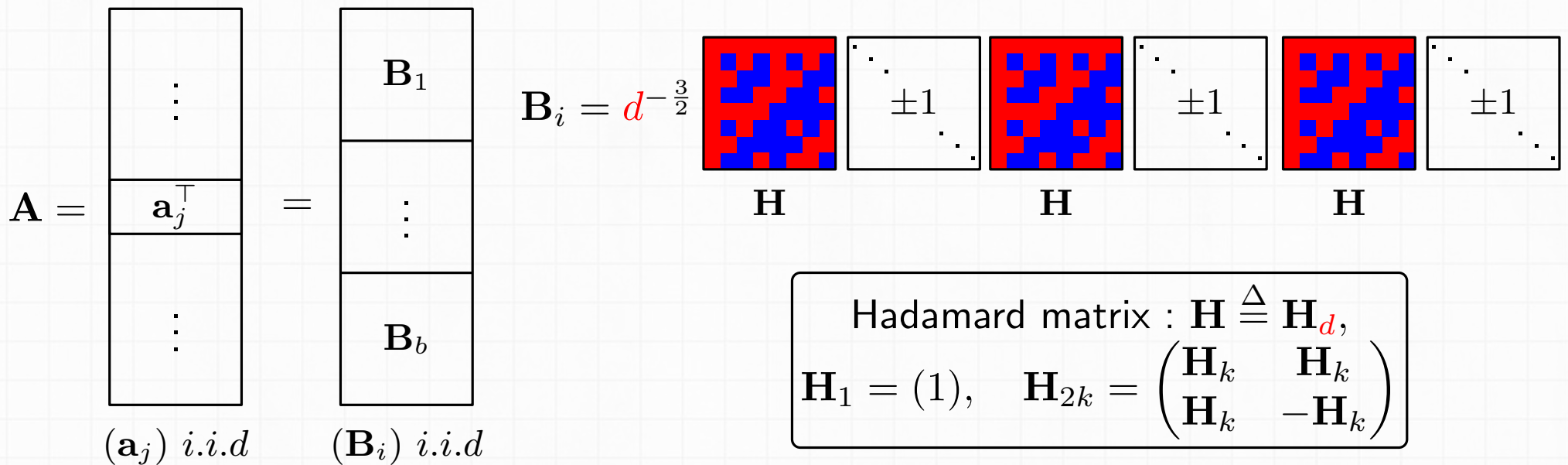
(Unsolvable)

→ Iterative algorithm
(inspired by *Proximal Gradient Descent*)



In practice (Encoding)

Reformulate : $\Phi(\mathbf{x}) = \frac{1}{m}(\mathbf{A}\mathbf{x})^{\odot 2}$

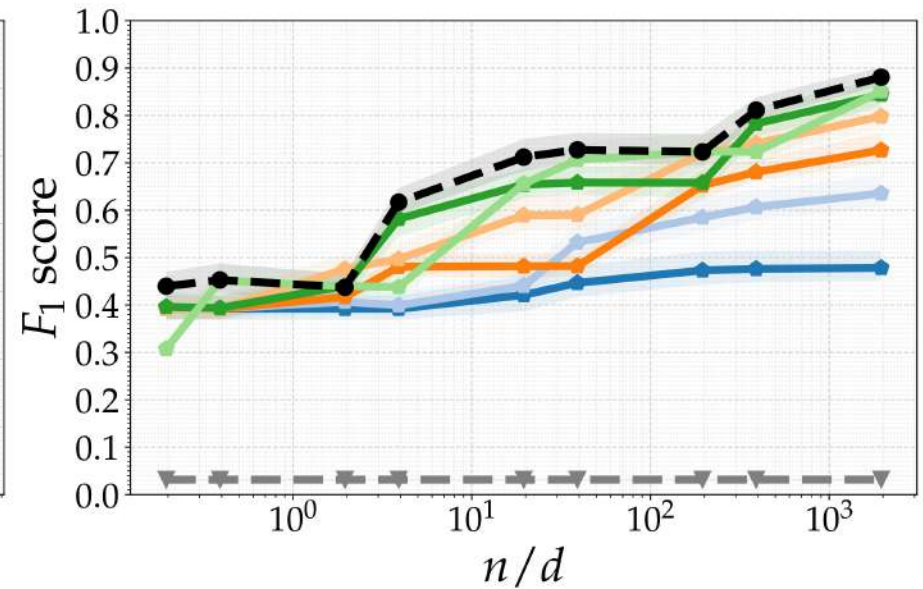
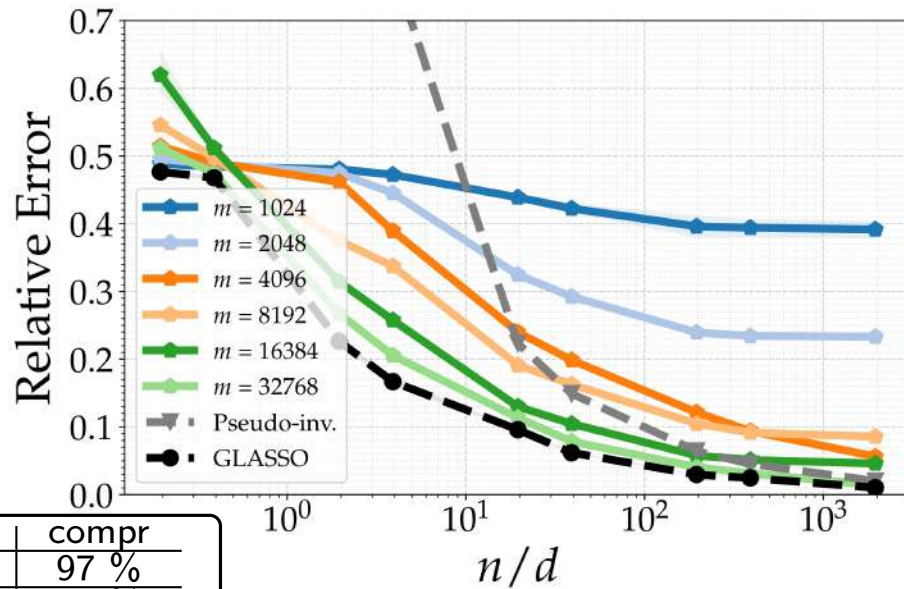


$\mathbf{H}\mathbf{x}$ can be computed in $\mathcal{O}(d \log d)$, without storing \mathbf{H} (fast Hadamard Transform).

	Comp.	Memory
Structured \mathbf{A}	$\mathcal{O}(m \log d)$	$\mathcal{O}(m)$

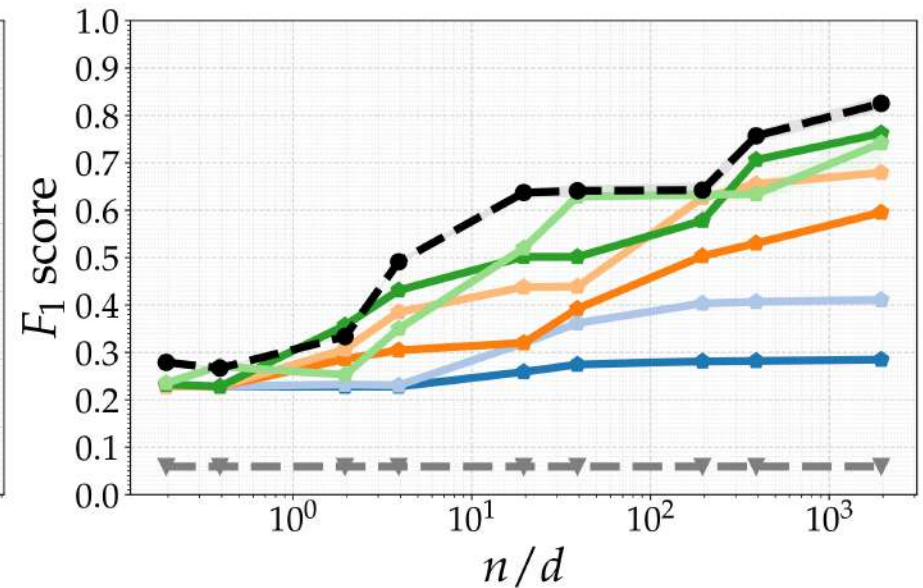
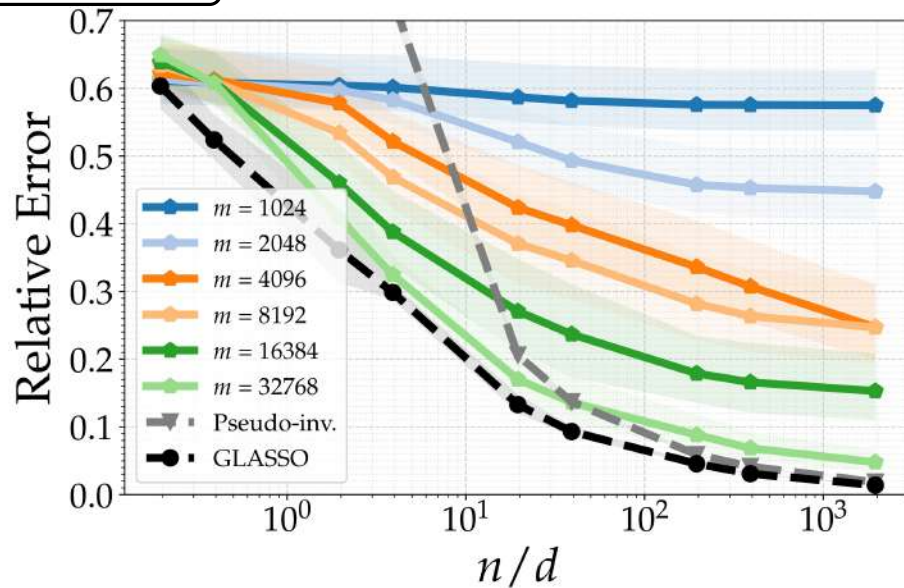
Experiments

Erdos ($d = 256$)



m	compr
1024	97 %
2048	94 %
4096	88 %
8192	75 %
16384	50 %

PowerLaw ($d = 256$)



Conclusion

Perspectives:

- Extending the theory to *structured* sketching operators.
- Guarantees for the practical decoder? (convergence? to what?)
- A more efficient practical decoder?
- Can we recover only properties of the graphs (e.g., clusters)

Preprint:

Compressive Recovery of Sparse Precision Matrices,
Vayer, L., Gribonval, Gonçalves (2023)
([arXiv:2311.04673](https://arxiv.org/abs/2311.04673))

Conclusion

Perspectives:

- Extending the theory to *structured* sketching operators.
- Guarantees for the practical decoder? (convergence? to what?)
- A more efficient practical decoder?
- Can we recover only properties of the graphs (e.g., clusters)

Preprint:

Compressive Recovery of Sparse Precision Matrices,
Vayer, L., Gribonval, Gonçalves (2023)
([arXiv:2311.04673](https://arxiv.org/abs/2311.04673))

Thank you for your attention.

In practice (Decoding)

A more practical decoding scheme:

$$f(\Sigma) \triangleq \frac{1}{2} \|\mathcal{A}(\Sigma) - \mathbf{s}\|_2^2 \text{ (sketch fidelity)}$$

$$\Sigma_0 \succ 0$$

$$\Sigma_{t+\frac{1}{2}} = \Sigma_t - \gamma \nabla f(\Sigma_t) \quad \text{(gradient descent step)}$$

$$\Sigma_{t+1} = \text{GLASSO}_{\gamma\lambda}[\Sigma_{t+\frac{1}{2}}] \quad \text{(denoising)}$$

In practice (Decoding)

A more practical decoding scheme:

$$f(\Sigma) \triangleq \frac{1}{2} \|\mathcal{A}(\Sigma) - \mathbf{s}\|_2^2 \text{ (sketch fidelity)}$$

$$\Sigma_0 \succ 0$$

$$\Sigma_{t+\frac{1}{2}} = \Sigma_t - \gamma \nabla f(\Sigma_t) \quad \text{(gradient descent step)}$$

$$\Sigma_{t+1} = \text{GLASSO}_{\gamma\lambda}[\Sigma_{t+\frac{1}{2}}] \quad \text{(denoising)}$$

Link with (Bregman) Proximal Gradient descent:

$$\arg \min_x \{f(x) + g(x)\}$$

init. x_0

$$x_{t+\frac{1}{2}} = x_t - \gamma \nabla f(x_t) \quad \text{(gradient descent step)}$$

$$x_{t+1} = \text{prox}_g(x_{t+\frac{1}{2}}) \quad \text{(proximal step)}$$

$$\curvearrowright \triangleq \arg \min_u \left\{ g(u) + \frac{1}{2} \|u - x_{t+\frac{1}{2}}\|_2^2 \right\}$$

In practice (Decoding)

A more practical decoding scheme:

$$f(\Sigma) \triangleq \frac{1}{2} \|\mathcal{A}(\Sigma) - \mathbf{s}\|_2^2 \text{ (sketch fidelity)}$$

$$\Sigma_0 \succ 0$$

$$\Sigma_{t+\frac{1}{2}} = \Sigma_t - \gamma \nabla f(\Sigma_t) \quad \text{(gradient descent step)}$$

$$\Sigma_{t+1} = \text{GLASSO}_{\gamma\lambda}[\Sigma_{t+\frac{1}{2}}] \quad \text{(denoising)}$$

Link with (Bregman) Proximal Gradient descent:

$$\arg \min_x \{f(x) + g(x)\}$$

init. x_0

$$x_{t+\frac{1}{2}} = x_t - \gamma \nabla f(x_t) \quad \text{(gradient descent step)}$$

$$x_{t+1} = \text{prox}_g^h(x_{t+\frac{1}{2}}) \quad \text{(Bregman proximal step)}$$

$$\triangleq \arg \min_u \left\{ g(u) + D_h(u | x_{t+\frac{1}{2}}) \right\}$$

$$\text{GLASSO}_{\lambda}[\mathbf{Z}] \triangleq \arg \min_{\Sigma \succ 0} \left\{ \lambda \|\Sigma^{-1}\|_{1,\text{off}} + D_h(\mathbf{Z} | \Sigma) \right\} \quad \text{with } h(\mathbf{X}) = -\log \det \mathbf{X}$$

Recipe of the Proof

Rewrite the RIP:

$$\text{RIP}(\delta, \mathfrak{S}_{k,a,b}) \Leftrightarrow \left| \|\mathcal{A}(\mathbf{U})\|_1 - 1 \right| \leq \delta, \quad \forall \mathbf{U} \in S[\mathfrak{S}_{k,a,b}]$$

$$\text{Normalized secant} : S[\mathfrak{S}_{k,a,b}] \triangleq \left\{ \frac{\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2}{\|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2\|}, \boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \mathfrak{S}_{k,a,b} \right\}$$

Goal: $\mathbb{P} \left(\left| \|\mathcal{A}(\mathbf{U})\|_1 - 1 \right| \leq \delta, \quad \forall \mathbf{U} \in S[\mathfrak{S}_{k,a,b}] \right) \geq 1 - \rho$

Ingredients:

- Pointwise concentration :

$$\forall \mathbf{U} \in S[\mathfrak{S}_{k,a,b}], \quad \forall t > 0, \quad \mathbb{P} \left(\left| \|\mathcal{A}(\mathbf{U})\|_1 - 1 \right| > t \right) \leq C(t)$$

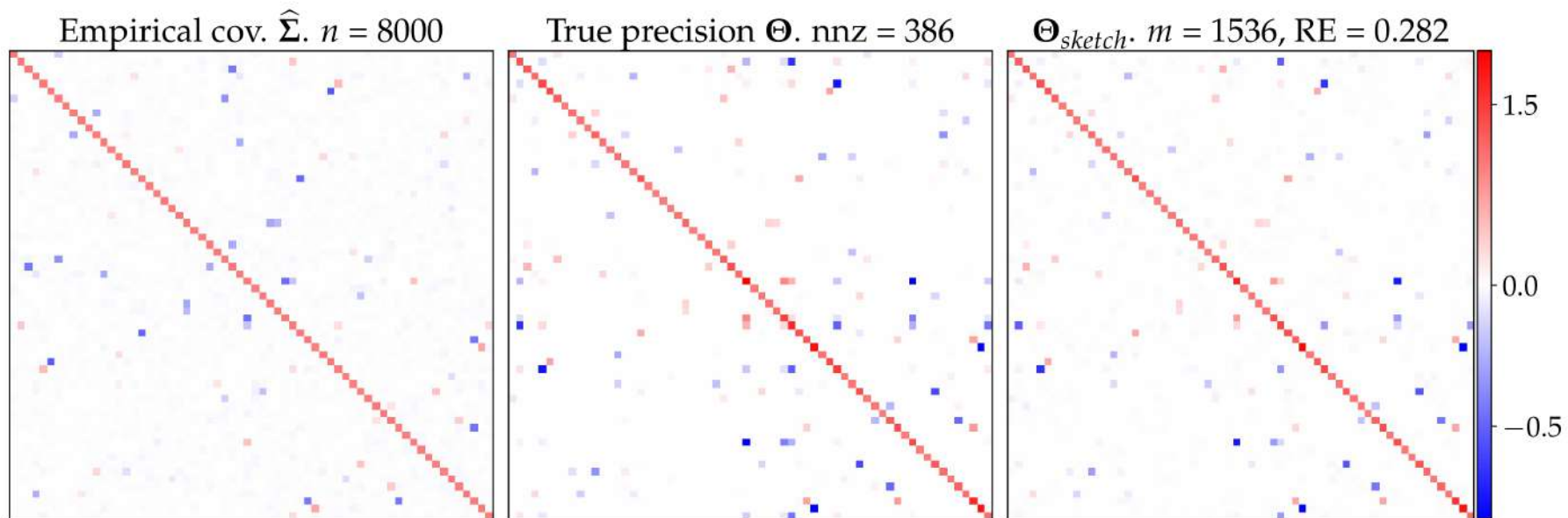
- Control of covering numbers:

$$\mathcal{N}(S[\mathfrak{S}_{k,a,b}], \|\cdot\|_{\Lambda}, \varepsilon) \triangleq \# \text{balls of size } \varepsilon \text{ to cover } S[\mathfrak{S}_{k,a,b}]$$

Experiments

$d = 64$, $m = 1536$ (compression 25%)

Erdos



PowerLaw

